

03

AUTOMATIC SOFTWARE DIVERSIFICATION FOR WEBASSEMBLY

Generating WebAssembly binaries is a process that commences with the original source code. This code is processed by a compiler, which can be conceptually divided into three primary components: a front-end that translates the source code into an intermediate representation, an optimizer that transforms this representation for the sake of performance, and a back-end that assembles the final WebAssembly binary. This modular architecture is perfectly implemented by LLVM, a compiler framework that supports a wide array of programming languages, including C, C++, and Rust.

Software Diversification, a proactive security measure, can be integrated at different points along this compilation pipeline. Yet, this strategy is limited if applied to the front-end. Specifically, it would need the development of a unique diversification mechanism for each language that can be compiled into WebAssembly, rendering it impractical for widespread use. For instance, in the case of LLVM, one would need to modify each of the front-ends to introduce diversification. Yet, implementing diversification at later stages of the compiler, such as the optimizer or the back-end, presents a more viable alternative. This is especially important given that approximately 70% of WebAssembly binaries are produced using LLVM-based compilers, as indicated by Hilbig et al. [?]. Therefore, featuring diversification techniques in the latter compiler stages is viable, leading us to propose compiler-based approaches to diversify WebAssembly binaries.

Our compiler-based strategies are highlighted in red and green in Figure 3.1. This approach introduces a diversifier component in the LLVM pipeline, generating LLVM IR variants and producing artificial software diversity for Wasm. This strategy encompasses two tools: CROW [?], which creates Wasm program variants, and MEWE [?], which merges these variants to provide multivariant execution [?] for Wasm.

On the other hand, diversification can be directly applied to the WebAssembly binary. This offers the most flexible approach since it is agnostic to the original programming language or the compiler employed. We propose a binary-based

⁰Comp. time 2023/10/02 09:45:41

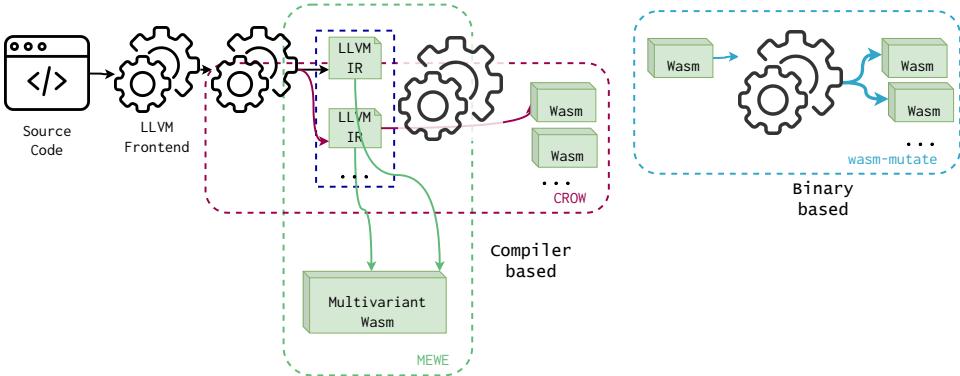


Figure 3.1: Approach landscape containing our three technical contributions: CROW squared in red, MEWE squared in green and WASM-MUTATE squared in blue. We annotate where our contributions, compiler-based and binary-based, stand in the landscape of generating WebAssembly programs.

approach, WASM-MUTATE, which is illustrated in blue in Figure 3.1. WASM-MUTATE [?] generates a pool of WebAssembly program variants through rewriting rules upon an e-graph [?] data structure. **TODO** Mention fingerprinting here !

To sum up, this dissertation contributes to the field of Software Diversification for WebAssembly through two strategies: compiler-based and binary-based. We present three main technical contributions: CROW, MEWE, and WASM-MUTATE, dissected upon in the subsequent sections. Moreover, we compare our technical contributions between them, showcasing that they are complementary and can be used in conjunction to provide a more robust diversification strategy. Furthermore, we outline the artifacts for our three technical contributions for the sake of open-research and reproducibility.

■ 3.1 CROW: Code Randomization of WebAssembly

This section details CROW [?], represented as the red squared tooling in Figure 3.1. CROW is designed to produce functionally equivalent Wasm variants from the output of an LLVM front-end, utilizing a custom Wasm LLVM backend.

Figure 3.2 illustrates CROW's workflow in generating program variants, a process compound of two core stages: *exploration* and *combining*. During the *exploration* stage, CROW processes every instruction within each function of the LLVM input, creating a set of functionally equivalent code variants. This process ensures a rich pool of options for the subsequent stage. In the *combining*

stage, these alternatives are assembled to form diverse LLVM IR variants, a task achieved through the exhaustive traversal of the power set of all potential combinations of code replacements. The final step involves the custom Wasm LLVM backend, which compiles the crafted LLVM IR variants into Wasm binaries.

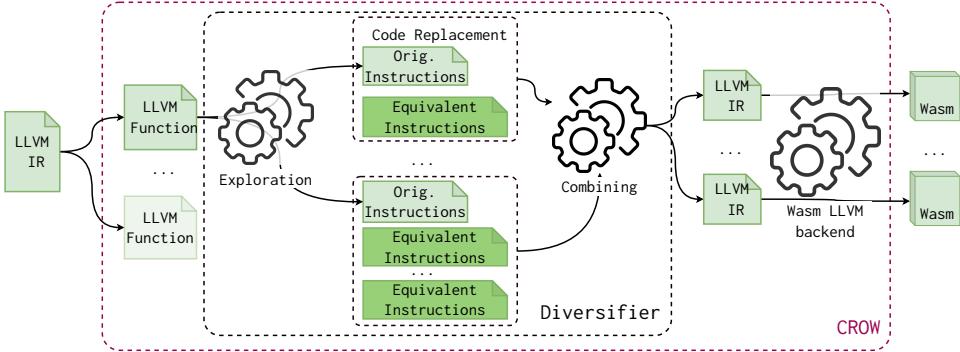


Figure 3.2: CROW components following the diagram in Figure 3.1. CROW takes LLVM IR to generate functionally equivalent code replacements. Then, CROW assembles program variants by combining them. Figure taken from [?].

■ 3.1.2 Variants’ space exploration

The cornerstone of CROW’s exploration mechanism is its code replacement generation strategy, which is inspired by the superdiversifier methodology proposed by Jacob et al. [?]. The search space for generating variants is delineated through an enumerative synthesis process, which systematically produces all possible code replacements for each instruction and its data flow graph in the original program. If a code replacement is identified to perform identically to the original program, it is flagged as a functionally equivalent variant. This equivalence is confirmed using a theorem solver for rigorous verification.

CROW is developed by modifying the enumerative synthesis implementation found in Souper [?], an LLVM-based superoptimizer. Specifically, CROW constructs a Data Flow Graph (DFG) for each LLVM instruction that returns an integer. Subsequently, it generates all viable expressions derived from a selected subset of the LLVM Intermediate Representation (IR) language for each DFG. The enumerative synthesis process incrementally generates code replacements, starting with the simplest expressions—those composed of a single instruction—and gradually increasing in complexity. The exploration process continues either until a timeout occurs or the size of the generated replacements exceeds a predefined threshold.

Notice that the searching space increases exponentially with the size of the LLVM IR language subset. Thus, we prevent CROW from synthesizing instructions without correspondence in the Wasm backend, effectively reducing the searching space. For example, creating an expression having the `freeze` LLVM instructions will increase the searching space for instruction without a Wasm's opcode in the end. Moreover, we disable the majority of the pruning strategies of Souper for the sake of more program variants. For example, Souper prevents the generation of commutative operations during the searching. On the contrary, CROW still uses such transformation as a strategy to generate program variants.

To generate variants, CROW employs a strategy of recombining all found code replacements. Leveraging the ascending nature of its enumerative synthesis process, CROW is capable of creating variants that may outperform the original program in both size and efficiency. For instance, the first functionally equivalent transformation identified is typically the most optimal in terms of code size. This approach offers developers a range of performance options, allowing them to balance between diversification and performance without compromising the latter. CROW applies these code transformations incrementally. For instance, if a suitable replacement is identified that can be applied at N different locations in the original program, CROW will generate N distinct program variants, each with the transformation applied at a unique location. This approach leads to a combinatorial explosion in the number of available program variants, especially as the number of possible replacements increases.

The last stage at CROW invokes a custom Wasm LLVM backend, which generates the Wasm programs. For it, we remove all built-in optimizations in the LLVM backend that could reverse Wasm variants, i.e., we disable all optimizations in the Wasm backend that could reverse the CROW transformations.

■ 3.1.3 Constant inferring

CROW, inherently adds a new transformation strategy that leads to more Wasm program variants, *constant inferring*. In concrete, Souper infers pieces of code as a single constant assignment and this is ported to CROW. This strategy mostly focuses on variables that are used to control branches. After a *constant inferring*, the generated program is considerably different from the original program, being suitable for diversification.

Let us illustrate the case with an example. The Babbage problem code in Listing 3.1 is composed of a loop that stops when it discovers the smallest number that fits with the Babbage condition in Line 4.

```

1 int babbage() {
2     int current = 0,
3         square;
4     while ((square=current*current) %4
5           ↪ 1000000 != 269696) {
6         current++;
7     printf ("The number is %d\n",
8           ↪ current);
9     return 0 ;
10 }
11 }
```

Listing 3.1: Babbage problem.
Taken from [?].

```

1 int babbage() {
2     int current = 25264;
3
4     printf ("The number is %d\n", current);
5     return 0 ;
6 }
```

Listing 3.2: Constant inferring transformation over the original Babbage problem in Listing 3.1. Taken from [?].

In theory, this value can also be inferred by unrolling the loop the correct number of times with the LLVM toolchain. However, standard LLVM tools cannot unroll the **while**-loop because the loop count is too large. CROW deals with this case, generating the program in Listing 3.2. It infers the value of **current** in Line 2 such that the Babbage condition is reached. Therefore, the condition in the loop will always be false. Then, the loop is dead code and is removed in the final compilation. The new program in Listing 3.2 is remarkably smaller and faster than the original code. Therefore, it offers differences both statically and at runtime¹

■ 3.1.4 CROW instantiation

Let us illustrate how CROW works with the example code in Listing 3.3. The **f** function calculates the value of $2 * x + x$ where **x** is the input for the function. CROW compiles this source code and generates the intermediate LLVM bitcode in the left most part of Listing 3.4. CROW potentially finds two integer returning instructions to look for variants, as the right-most part of Listing 3.4 shows.

```

1 int f(int x) {
2     return 2 * x + x;
3 }
```

Listing 3.3: C function that calculates the quantity $2x + x$.

¹Notice that for the sake of illustration, we show both codes in C language, this process inside CROW is performed directly in LLVM IR.

```

define i32 @f(i32) {           Replacement candidates      Replacement candidates for
                                for code_1                code_2
    %2 = mul nsw i32 %0,2
    %3 = add nsw i32 %0,%2    %2 = mul nsw i32 %0,2      %3 = add nsw i32 %0,%2
                                %2 = add nsw i32 %0,%0    %3 = mul nsw %0, 3:i32
                                ret i32 %3
                                %2 = shl nsw i32 %0, 1:i32
}
define i32 @main() {
    %1 = tail call i32 @f(
        i32 10)
    ret i32 %1
}

```

Listing 3.4: LLVM’s intermediate representation program, its extracted instructions and replacement candidates. Gray highlighted lines represent original code, green for code replacements.

```

%2 = mul nsw i32 %0,2          %2 = mul nsw i32 %0,2
%3 = add nsw i32 %0,%2         %3 = mul nsw %0, 3:i32
                                %2 = add nsw i32 %0,%0
                                %3 = mul nsw %0, 3:i32
                                %2 = shl nsw i32 %0, 1:i32
                                %3 = add nsw i32 %0,%2
                                %2 = shl nsw i32 %0, 1:i32
                                %3 = mul nsw %0, 3:i32

```

Listing 3.5: Candidate code replacements combination. Orange highlighted code illustrate replacement candidate overlapping.

CROW, detects `code_1` and `code_2` as the enclosing boxes in the left most part of Listing 3.4 shows. CROW synthesizes $2 + 1$ candidate code replacements for each code respectively as the green highlighted lines show in the right most parts of Listing 3.4. The baseline strategy of CROW is to generate variants out of all possible combinations of the candidate code replacements, *i.e.*, uses the power set of all candidate code replacements.

In the example, the power set is the cartesian product of the found candidate code replacements for each code block, including the original ones, as Listing 3.5 shows. The power set size results in 6 potential function variants. Yet, the generation stage would eventually generate 4 variants from the original program. CROW generated 4 statically different Wasm files, as Listing 3.6 illustrates. This gap between the potential and the actual number of variants is a consequence of the redundancy among the bitcode variants when composed into one. In other words, if the replaced code removes other code blocks, all possible combinations having it will be in the end the same program. In the example case, replacing `code_2` by `mul nsw %0, 3`, turns `code_1` into dead code, thus, later replacements generate the same program variants. The rightmost part of Listing 3.5 illustrates how for three different combinations, CROW produces the same variant. We call this phenomenon a *code replacement overlapping*.

```

func $f (param i32) (result i32)
    local.get 0
    i32.const 2
    i32.mul
    local.get 0
    i32.add

func $f (param i32) (result i32)
    local.get 0
    local.get 0
    i32.add
    local.get 0
    i32.add

func $f (param i32) (result i32)
    local.get 0
    i32.const 1
    i32.shl
    local.get 0
    i32.add

func $f (param i32) (result i32)
    local.get 0
    i32.const 3
    i32.mul

```

Listing 3.6: Wasm program variants generated from program Listing 3.3.

One might think that a reasonable heuristic could be implemented to avoid such overlapping cases. Instead, we have found it easier and faster to generate the variants with the combination of the replacement and check their uniqueness after the program variant is compiled. This prevents us from having an expensive checking for overlapping inside the CROW code. Still, this phenomenon calls for later optimizations in future works.

Contribution paper and artifact

CROW fully presented in Cabrera-Arteaga et al. "CROW: Code Randomization of WebAssembly" *Network and Distributed System Security Symposium, MADWeb* <https://doi.org/10.14722/madweb.2021.23004>.

CROW source code is available at <https://github.com/ASSERT-KTH/slumps>

■ 3.2 MEWE: Multi-variant Execution for WebAssembly

This section describes MEWE [?]. MEWE synthesizes diversified function variants by using CROW. It then provides execution-path randomization in a Multivariant Execution (MVE) [?]. MEWE generates application-level multivariant binaries without changing the operating system or Wasm runtime. It creates an MVE by intermixing functions for which CROW generates variants, as illustrated by the green square in Figure 3.1. MEWE inlines function variants when appropriate, resulting in call stack diversification at runtime.

As illustrated in Figure 3.3, MEWE takes the LLVM IR variants generated by CROW's diversifier. It then merges LLVM IR variants into a Wasm multivariant. In the figure, we highlight the two components of MEWE,

Multivariant Generation and the *Mixer*. In the *Multivariant Generation* process, MEWE gathers the LLVM IR variants created by CROW.

The *Mixer* has three specific steps: 1) link the LLVM multivariant binary, 2) inject a random generator and 3) tamper the application’s entrypoint. The Mixer component of MEWE links the multivariant binary and creates a new entrypoint for the binary called *entrypoint tampering*. The tampering is needed in case the output of CROW are variants of the original entrypoint, e.g. the *main* function. Concretely, it wraps the dispatcher for the entrypoint variants as a new function for the final Wasm binary and is declared as the application entrypoint. The random generator is needed to perform the execution-path randomization. For the random generator, we rely on WASI’s specification [?] for the random behavior of the dispatchers. However, its exact implementation is dependent on the platform on which the binary is deployed. Finally, using the same custom Wasm LLVM backend as CROW, we generate a standalone multivariant Wasm binary. Once generated, the multivariant Wasm binary can be deployed to any Wasm engine.

■ 3.2.2 Multivariant generation

The key component of MEWE consists of combining the variants into a single binary. The goal is to support execution-path randomization at runtime [?]. The core idea is to introduce one dispatcher function per original function with variants. A dispatcher function is a synthetic function in charge of choosing a variant at random when the original function is called. With the introduction of the dispatcher function, MEWE turns the original call graph into a multivariant call graph, defined as follows.

In Figure 3.4, we show the original static call graph for an original program (top of the figure), as well as the multivariant call graph generated with MEWE (bottom of the figure). The gray nodes represent function variants, the green nodes function dispatchers, and the yellow nodes are the original functions. The directed edges represent the possible calls. The original program includes three functions. MEWE generates 43 variants for the first function, none for the second, and three for the third. MEWE introduces two dispatcher nodes for the first and third functions. Each dispatcher is connected to the corresponding function variants to invoke one variant randomly at runtime.

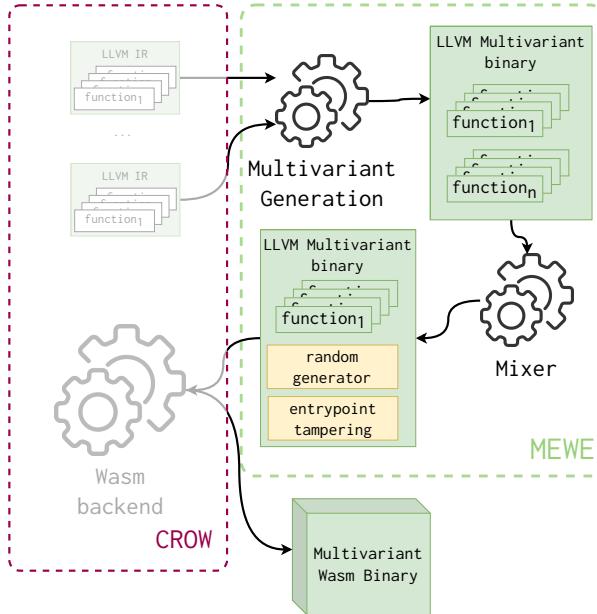


Figure 3.3: Overview of MEWE workflow. It takes as input an LLVM binary. It first generates a set of functionally equivalent variants for each function in the binary using CROW. Then, MEWE generates an LLVM multivariant binary composed of all the function variants. Finally, the Mixer includes the behavior in charge of selecting a variant when a function is invoked. Finally, the MEWE mixer composes the LLVM multivariant binary with a random number generation library and tampers the original application entrypoint. The final process produces a Wasm multivariant binary ready to be deployed. Figure partially taken from [?].

```

2 ; Multivariant foo wrapping ;
3 define internal i32 @foo(i32 %0) {
4   entry:
5     ; It first calls the dispatcher to discriminate between the created
      variants ;
6     %1 = call i32 @discriminate(i32 3)
7     switch i32 %1, label %end [
8       i32 0, label %case_43_
9       i32 1, label %case_44_
10    ]
11   ;One case for each generated variant of foo ;
12   case_43_:
13     %2 = call i32 @foo_43_(%0)
14     ret i32 %2
15   case_44_:
16     ; MEWE can inline the body of the a function variant ;
17     %3 = <body of foo_44_ inlined>
18     ret i32 %3
19 end:
20   ; The original is also included ;
21   %4 = call i32 @foo_original(%0)
22   ret i32 %4
23 }
```

Listing 3.7: Dispatcher function embedded in the multivariant binary of the

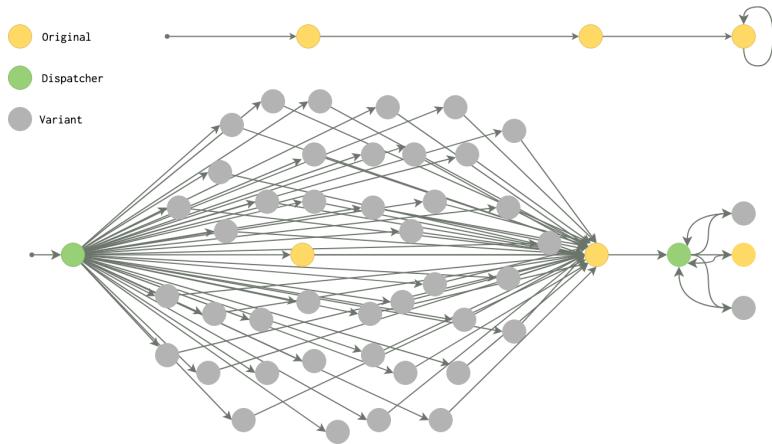


Figure 3.4: Example of two static call graphs. At the top, is the original call graph, and at the bottom, is the multivariant call graph, which includes nodes that represent function variants (in gray), dispatchers (in green), and original functions (in yellow). Figure taken from [?].

TODO Recheck these works on the decision of using switch cases.

In Listing 3.7, we demonstrate how MEWE constructs the function dispatcher, corresponding to the rightmost green node in Figure 3.4, which handles three created variants including the original. The dispatcher function retains the same signature as the original function. Initially, the dispatcher invokes a random number generator, the output of which is used to select a specific function variant for execution (as seen on line 6 in Listing 3.7). To enhance security, we employ a switch-case structure within the dispatcher, mitigating vulnerabilities associated with speculative execution-based attacks [?] (refer to lines 12 to 19 in Listing 3.7). This approach also eliminates the need for multiple function definitions with identical signatures, thereby reducing the potential attack surface in cases where the function signature itself is vulnerable [?]. Additionally, MEWE can inline function variants directly into the dispatcher, obviating the need for redundant definitions (as illustrated on line 16 in Listing 3.7). Remarkably, we prioritize security over performance, i.e., while using indirect calls in place of a switch-case could offer constant-time performance benefits, we implement switch-case structures.

In Listing 3.7, we illustrate the LLVM construction for the function dispatcher corresponding to the right most green node of Figure 3.4. Notice that, the dispatcher function is constructed using the same signature as the original function. It first calls the random generator, which returns a value used to invoke a specific function variant (see line 6 in Listing 3.7). We utilize a switch-case structure in the dispatchers to prevent indirect calls, which are vulnerable to speculative execution-based attacks [?] (see lines 12 to 19 in Listing 3.7), i.e.,

the choice of a switch-case also avoids having multiple function definitions with the same signature, which could increase the attack surface in case the function signature is vulnerable [?]. In addition, MEWE can inline function variants inside the dispatcher instead of defining them again (see line 16 in Listing 3.7). Remarkably, we trade security over performance since dispatcher functions that perform indirect calls, instead of a switch-case, could improve the performance of the dispatchers as indirect calls have constant time.

Contribution paper and artifact

MEWE is fully presented in Cabrera-Arteaga et al. "Multi-Variant Execution at the Edge" *Proceedings of ACM, Moving Target Defense* <https://dl.acm.org/doi/abs/10.1145/3560828.3564007>

MEWE is also available as an open-source tool at <https://github.com/ASSERT-KTH/MEWE>

■ 3.3 WASM-MUTATE: Fast and Effective Binary for WebAssembly

In this section, we introduce our third technical contribution, WASM-MUTATE [?], a tool that generates functionally equivalent variants of a WebAssembly binary input. Leveraging rewriting rules and e-graphs [?] for diversification space traversals, WASM-MUTATE synthesizes program variants by transforming parts of the original binary. In Figure 3.1, we highlight WASM-MUTATE as the blue squared tooling for a visual representation.

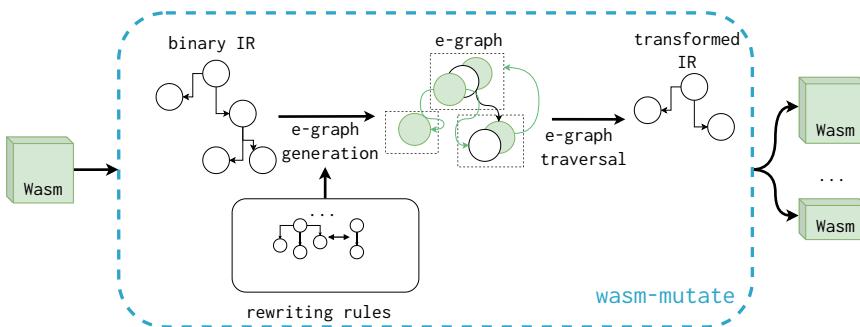


Figure 3.5: WASM-MUTATE high-level architecture. It generates functionally equivalent variants from a given WebAssembly binary input. Its central approach involves synthesizing these variants by substituting parts of the original binary using rewriting rules, boosted by diversification space traversals using e-graphs.

Figure 3.5 illustrates the workflow of WASM-MUTATE, which initiates with a WebAssembly binary as its input. The first step involves parsing this binary to create suitable abstractions, e.g. an intermediate representation. Subsequently, WASM-MUTATE utilizes predefined rewriting rules to construct an e-graph for the initial program, encapsulating all potential equivalent codes derived from the rewriting rules. Then, pieces of the original program are randomly substituted by the result of random e-graph traversals, resulting in a variant that maintains functional equivalence to the original binary. This assurance of functional preservation is rooted in the inherent properties of the individual rewrite rules employed.

■ 3.3.2 WebAssembly Rewriting Rules

The WASM-MUTATE framework incorporates a comprehensive set of 135 rewriting rules, organized into distinct categories as meta-rules. These rewriting rules are inspired by the foundational work of Sasnauskas et al. [?], and are extended to include a predicate that enforces the conditions under which a replacement can occur. Each rule is structured as a tuple, denoted as $(\text{LHS}, \text{RHS}, \text{Cond})$. Here, **LHS** specifies the segment of code targeted for replacement, **RHS** describes its functionally equivalent substitute, and **Cond** outlines the conditions that must be met for the replacement to take place. In the context of WebAssembly binaries, the **Cond** predicate ensures that any replacement adheres to type constraints.

WASM-MUTATE features seven different meta-rules, ranging from high-level changes affecting all binary section structure to low-level modifications within the code section, including alterations to both control and data flow. To simplify the discussion, the subsequent text will focus on one specific meta-rule implemented in WASM-MUTATE, the **Peephole** meta-rule. For an in-depth explanation of the remaining meta-rules, refer to [?].

The **Peephole** meta-rule focuses on the rewriting of instruction sequences found within function bodies, representing the lowest level of rewriting. In WASM-MUTATE, we have devised 125 rewriting rules specifically for this category. WASM-MUTATE is structured to ensure the determinism of the instructions selected for replacement. Therefore, any rewriting rule inside the Peephole meta-rule avoids instructions that might induce undefined behavior, e.g., function calls. Consequently, preserving the original functionality of the control frame labels (see Subsection 2.1.3).

Moreover, we augment the internal representation of a Wasm program to bolster WASM-MUTATE’s transformation capabilities through the **Peephole** meta-rule. Concretely, we augment the parsing stage in WASM-MUTATE by including custom operator instructions. These custom operator instructions are designed to bolster WASM-MUTATE by utilizing well-established code diversification techniques through rewriting rules. In practice, we add 4 custom operator instructions, namely **container**, **useglobal**, **unfold**, and **rand**. Custom

operands highlight the versatility of WASM-MUTATE as a general-purpose binary rewriting engine.

In the example below, we demonstrate a rewriting rule of the Peephole meta-rule that leverages a custom operator to insert `nop` instructions into any WebAssembly program place, a well-known low-level diversification strategy [?], while using the `container` custom operator:

LHS `x`

RHS (`container (x nop)`)

In practice, custom operators are only part of the rewriting rules of WASM-MUTATE. This means that, when converting Wasm to the intermediate representation no custom operator is generated. When converting back to the WebAssembly binary format from the intermediate representation, custom instructions are meticulously handled to retain the original functionality of the WebAssembly program. For example, the `container` custom operator is removed while its operands are encoded back to Wasm in their corresponding opcodes.

■ 3.3.3 E-Graphs traversals

We developed WASM-MUTATE leveraging e-graphs, a specific graph data structure for representing rewriting rules [?]. In the context of WASM-MUTATE, the e-graph is constructed from a WebAssembly program. This entails the transformation of each distinct expression, operator, and operand into e-nodes. A primer e-graph is built from the original program. This initial e-graph is subsequently augmented with e-nodes and e-classes derived from each one of the rewriting rules (we detailed the e-graph construction process in Section 3 of [?]).

Willsey et al. highlight the potential for high flexibility in extracting code fragments from e-graphs, a process that can be recursively orchestrated through a cost function applied to e-nodes and their respective operands. This methodology ensures the functional equivalence of the derived code [?]. For instance, e-graphs solve the problem of providing the best code out of several optimization rules [?]. To extract the "optimal" code from an e-graph, one might commence the extraction at a specific e-node, subsequently selecting the AST with the minimal size from the available options within the corresponding e-class's operands. In omitting the cost function from the extraction strategy leads us to a significant property: *any path navigated through the e-graph yields a functionally equivalent code variant.*

We exploit such property to fastly generate diverse WebAssembly variants. We propose and implement an algorithm that facilitates the random traversal of an e-graph to yield functionally equivalent program variants, as detailed in Algorithm 1. This algorithm operates by taking an e-graph, an e-class node (starting with the root's e-class), and a parameter specifying the maximum

extraction depth of the expression, to prevent infinite recursion. Within the algorithm, a random e-node is chosen from the e-class (as seen in lines 5 and 6), setting the stage for a recursive continuation with the offspring of the selected e-node (refer to line 8). Once the depth parameter reaches zero, the algorithm extracts the most concise expression available within the current e-class (line 3). Following this, the subexpressions are built (line 10) for each child node, culminating in the return of the complete expression (line 11).

Algorithm 1 e-graph traversal algorithm taken from [?].

```

1: procedure TRAVERSE(egraph, eclass, depth)
2:   if depth = 0 then
3:     return smallest_tree_from(egraph, eclass)
4:   else
5:     nodes  $\leftarrow$  egraph[eclass]
6:     node  $\leftarrow$  random_choice(nodes)
7:     expr  $\leftarrow$  (node, operands = [])
8:     for each child  $\in$  node.children do
9:       subexpr  $\leftarrow$  TRAVERSE(egraph, child, depth - 1)
10:      expr.operands  $\leftarrow$  expr.operands  $\cup$  {subexpr}
11:    return expr
```

■ 3.3.4 WASM-MUTATE instantiation

Let us illustrate how WASM-MUTATE generates variant programs by using the before enunciated algorithm. Here, we use Algorithm 1 with a maximum depth of 1. In Listing 3.8 a hypothetical original Wasm binary is illustrated. In this context, a potential user has set two pivotal rewriting rules: `(x, container (x nop),)` and `(x, x i32.add 0, x instanceof i32)`. The former rule, which has been previously discussed in Subsection 3.3.2, grants the ability to append a `nop` instruction to any subexpression. Conversely, the latter rule adds zero to any numeric value .

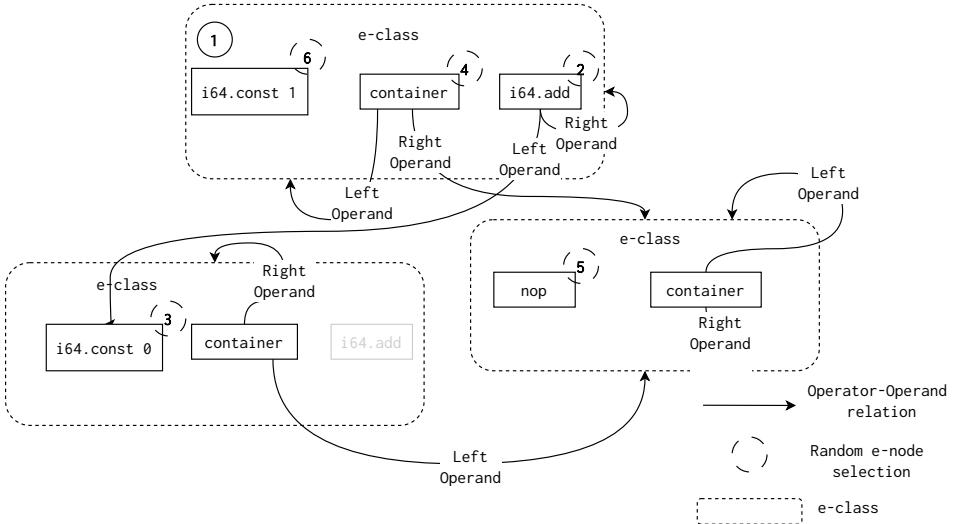


Figure 3.6: e-graph built for rewriting the first instruction of Listing 3.8.

```
(module
  (type (;0;) (func (param i32 f32) (result i64)))
  (func (;0;) (type 0) (param i32 f32) (result i64)
    i64.const 1)
)
```

Listing 3.8: Wasm function.

```
(module
  (type (;0;) (func (param i32 f32) (result i64)))
  (func (;0;) (type 0) (param i32 f32) (result i64)
    (i64.add (
      i64.const 0
      i64.const 1
      nop
    )))
)
```

Listing 3.9: Random peephole mutation using egraph traversal for Listing 3.8 over e-graph Figure 3.6. The textual format is folded for better understanding.

Leveraging the code presented in Listing 3.8 alongside the defined rewriting rules, we build the e-graph, simplified in Figure 3.6. In the figure, we highlight various stages of Algorithm 1 in the context of the scenario previously described. The algorithm initiates at the e-class with the instruction `i64.const 1`, as seen in Listing 3.8. At ②, it randomly selects an equivalent node within the e-class, in this instance taking the `i64.add` node, resulting: `expr`

= `i64.add 1 r`. As the traversal advances, it follows on the left operand of the previously chosen node, settling on the `i64.const 0` node within the same e-class ③. Then, the right operand of the `i64.add` node is chosen, selecting the `container` ④ operator yielding: `expr = i64.or (i64.const 0 container (r nop))`. The algorithm chooses the right operand of the `container` ⑤, which correlates to the initial instruction e-node highlighted in ⑥, culminating in the final expression: `expr = i64.or (i64.const 0 container(i64.const 1 nop)) i64.const 1`. As we proceed to the encoding phases, the `container` operator is ignored as a real Wasm instruction, finally resulting in the program in Listing 3.9.

Notice that, within the e-graph showcased in Figure 3.6, the `container` node maintains equivalence across all e-classes. Consequently, increasing the depth parameter in Algorithm 1 would potentially escalate the number of viable variants infinitely.

Contribution paper and artifact

WASM-MUTATE is fully presented in Cabrera-Arteaga et al. "WASM-MUTATE: Fast and Effective Binary Diversification for WebAssembly"
<https://arxiv.org/pdf/2309.07638.pdf>.

WASM-MUTATE is available at <https://github.com/bytocodealliance/wasm-tools/tree/main/crates/wasm-mutate> as a contribution to the bytocodealliance organization ^a.

^a<https://bytocodealliance.org/>

■ 3.4 Comparing CROW, MEWE, and WASM-MUTATE

In this section, we compare CROW, MEWE, and WASM-MUTATE, highlighting their key differences. These distinctions are summarized in Table 3.1. The table is organized into columns that represent attributes of each tool: the tool's name, input format, core diversification strategy, number of variants generated within an hour, targeted sections of the WebAssembly binary for diversification, strength of the generated variants, and the security applications of these variants. Each row in the table corresponds to a specific tool. Notice that, the data and insights presented in the table are sourced from the respective papers of each tool and, from the previous discussion in this chapter.

Tool	Input	Core	Variants in 1h	Target	Variants Strength	Security applications
CROW	Source code or LLVM Ir	Enumerative synth.	> 1k	Code section	96%	Resilience against: signature-based identification, static analysis and side-channel attacks.
MEWE	Source code or LL VM Ir	CROW + Execution path randomization	> 1k	Code + Function sections	96%	Resilience against: signature-based identification, static and dynamic analysis, web timing-based attacks.
WASM-MUTATE	rewriting rules + Wasm bin.	e-graph random traversals	> 10k	Any Wasm part	76%	Resilience against: signature-based identification, static analysis, fingerprinting and timing side-channel attacks.

Table 3.1: Comparing CROW, MEWE and WASM-MUTATE. The table columns are: the tool’s name, input format, core diversification strategy, number of variants generated within an hour, targeted sections of the WebAssembly binary, strength of the generated variants, and the security applications of these variants. The Variant strength accounts for the capability of each tool on generating variants that are preserved after the JIT compilation of V8 and wasmtime in average. Our three technical contributions are complementary tools that can be combined.

CROW is a compiler-based strategy, needing access to the source code or its LLVM IR representation to work. Its core is an enumerative synthesis implementation with functionality verification using SMT solvers, ensuring the functional equivalence of the generated variants. In addition, MEWE extends the capabilities of CROW, utilizing the same underlying technology to create program variants. It goes a step further by packaging the LLVM IR variants into a Wasm multivariant, providing MVE through execution path randomization. Both CROW and MEWE are fully automated, requiring no user intervention besides the input source code. WASM-MUTATE, on the other hand, is a semi-automated, binary-based tool. It needs a set of rewriting rules and the Wasm binary as inputs to generate program variants, centralizing its core around random e-graph traversals. Remarkably, WASM-MUTATE removes the need for compiler adjustments, offering compatibility with any existing WebAssembly binary.

WASM-MUTATE generates more unique variants in one hour than CROW and MEWE in at least one order of magnitude. This is because WASM-MUTATE is a binary-based tool, not requiring any compilation step further the lazy parsing of the diversification target section. Besides, WASM-MUTATE can generate variants in any part of the Wasm binary, while CROW and MEWE are limited to the code and function sections. In addition, CROW and MEWE generation capabilities are limited by the *overlapping* phenomenon discussed in Subsection 3.1.4. On the other hand, CROW and MEWE use enumerative synthesis and verify functional equivalence through SMT solvers. This approach not only has the potential to exceed handcrafted optimizations [?] but also ensures that the transformations are preserved. In other words, the transformations generated out of CROW and MEWE are virtually irreversible by JIT compilers, such as V8 and wasmtime. This phenomenon is highlighted in the *Variants strength* column of Table 3.1, where we show that CROW and MEWE generate variants with 96% of preservation against 75% of WASM-MUTATE.

■ 3.4.2 Security applications

The last column of Table 3.1 highlights the security applications of the variants generated by our three technical contributions. Our tools generate many different and highly preserved code variants. This means that these variants, each with unique WebAssembly codes, maintain their distinctiveness even after JIT compilers translate them into machine codes (see Subsection 2.1.6). The preservation feature significantly reduces the impact of side-channel attacks that exploit specific machine code instructions, e.g., port contention [?]. Besides, the preserved transformations of the generated variants serve to reduce potential attack surfaces, thereby impeding signature-based identification.

Altering the layout of a WebAssembly program inherently influences its managed memory during runtime, a component not overseen by the WebAssembly program itself (see Definition 1). This phenomenon is especially important for CROW and MEWE, given that they do not directly address the WebAssembly

memory model. Significantly, CROW and MEWE considerably alter the managed memory by modifying the layout of the WebAssembly program. For example, the *constant inferring* transformations significantly alter the layout of program variants, affecting unmanaged memory elements such as the returning address of a function. Furthermore, WASM-MUTATE not only affects managed memory through changes in the WebAssembly program layout. It also adds rewriting rules to transform unmanaged memory instructions. Memory alterations, either to the unmanaged or managed memories, have substantial security implications, by eliminating potential jump points that facilitate malicious activities within the binary [?].

Besides, our technical contributions enhance security against timing-based attacks by creating variants that exhibit a wide range of execution times. This strategy is especially prominent in MEWE’s approach, which develops multivariants functioning on randomizing execution paths, thereby thwarting attempts at timing-based inference attacks [?]. Adding another layer benefit, the integration of diverse variants into multivariants can potentially disrupt dynamic analysis tools such as symbolic executors [?]. Concretely, different control flows through a random discriminator, exponentially increase the number of possible execution paths, making multivariant binaries virtually unexplorable.

Key Takeaway

Our three technical contributions serve as complementary tools that can be combined to create a more comprehensive and robust software diversification strategy. For instance, when the source code for a WebAssembly binary is either non-existent or inaccessible, WASM-MUTATE offers a viable solution for generating code variants. On the other hand, CROW and MEWE excel in scenarios where high preservation is crucial, particularly when the generated variants may be subject to further analysis. Furthermore, WASM-MUTATE can benefit from the enumerative synthesis techniques employed by CROW and MEWE. Specifically, WASM-MUTATE could incorporate the transformations generated by these tools as rewriting rules.

■ 3.5 Conclusions

In this chapter, we discuss the technical specifics underlying our primary technical contributions. We elucidate the mechanisms through which CROW generates program variants. Subsequently, we discuss MEWE, offering a detailed examination of its role in forging MVE for WebAssembly. We also explore the details of WASM-MUTATE, highlighting its pioneering utilization of an e-graph traversal algorithm to spawn Wasm program variants. Remarkably, we undertake a comparative analysis of the three tools, highlighting their respective benefits

and limitations, alongside the potential security applications of the generated Wasm variants.

In ??, we present four use cases that support the exploitation of these tools. ?? serves to bridge theory with practice, showcasing the tangible impacts and benefits realized through the deployment of CROW, MEWE, and WASM-MUTATE.